

Final Project Milestone:
Real Time Tennis Match Prediction Using Machine Learning

Yang "Eddie" Chen, Yi Zhong, Yubo Tian

November 21, 2017

1. Motivation

Tennis matches are fun to watch because they are full of surprises. In the 2017 Stuttgart Open, Roger Federer, then an 18-time grand slam champion, lost to the world No. 302 player Tommy Haas. Federer lost the opening match at a grass-court tournament, a phenomenon that hadn't happened since 2002. How can we predict a rare loss like this? Based on past performance alone, any model would have predicted a Federer win in pre-game bets. This motivates us to apply machine learning to predict tennis matches in real time; in particular, we want to explore how well we can combine both historical performance data and real-time, just-happened set-by-set outcome, to achieve a better prediction. Anecdotally, most of us probably have been consumers of ESPN's in-game predictions when watching an NFL game on ESPN's website, with a win chance chart displaying the real-time probability that ESPN thinks each team has. We hope, with a similar hybrid model that combines historical and real-time information, both tennis players and sports bettors can benefit from better predictions and new insights.

Tennis is an ideal candidate for a hierarchical model as a match consists of a sequence of sets, which in turn consist of a sequence of games, which in turn consist of a sequence of points. Extensive researches have been done on tennis pre-game predictions. ATP provides prediction before game based on past performance (including previous encounter, current ranking, etc). Utilizing existing data and building on top of existing research such as Sipko [8], Clarke and Dyte [1] and Knottenbelt [4], this paper seeks to apply classification algorithms from machine learning to model men's professional singles matches innovatively by using both pre-game, historical performance calculated via a common opponent model [4], and real-time, on-court stats [3] [5]. As pointed out by Sipko [8], an in-game "approach to tennis match prediction can be more accurate, as it allows the

model to capture the change in a player’s performance over the course of the match. For example, different players fatigue during a match in different ways." We hope that results from our predictions can be extended to give real-time coaching advice to support game strategy decision.

2. Method

2.1. Feature Extraction

2.1.1. Feature Representation

A supervised machine learning algorithm requires a set of labeled examples for training. In our project, each training example combines both the historical aspect, which is the stats calculated from the past, and the real-time aspect, which is the immediately preceding set information. (Reminder that our project focuses on predicting match outcome from both the historical data and the real-time info). In terms of representation, the training example is made up of:

1. A vector of input features (X), describing the players and the match
2. A target value (y), indicating the match outcome

In any given singles match, we have two players in one game. We denote them as Player 1 and Player 2, and thus define the target value as $y = 1$ if Player 1 wins, 0 otherwise. In the datasets we found online, it is usually the case that Player 1 is denoted as the winner and Player 2 being the loser. In this case, we duplicated the entire dataset but swapped Player 1 and Player 2 and set $y = 0$, as well as inverted all difference features. This way, we still maintained our notation, but also made sure the model has both 1 and 0 labels to train on.

2.1.2. Symmetric Match & Set Feature Representation

As suggested by Sipko [8] and Clarke and Dyte [1], when considering the characteristics of *both* players in a given match, we are mainly interested in the *difference* between two players for our variables of interest. Take for example, the rank information between two players for a given match. We construct a single feature called $RANK_{DIFF} = RANK_1 - RANK_2$, where $RANK_1$ and $RANK_2$ are the ranks of players 1 and 2 respectively at the time of the match. We selected this design (as opposed to the alternative of including both players’ info, such as both $RANK_1$ and $RANK_2$), because previous research has shown that the difference variables are predictive enough [8] [2], and this way we cut down the feature size by half as compared to the alternative. Reduced feature column count reduces the *variance* of the model. Moreover, using difference variables allows us to have a symmetric model: it will produce the same results even if we swap Player 1 and Player 2, and define $y = 1$ if Player 2 wins, 0 otherwise. Should we have two different variables such as $RANK_1$ and $RANK_2$ for an asymmetric model, the model might assign different weights to the two variables despite them referring to the same

fundamental feature, only differing by the player labelling. Having a *symmetric* model eliminates the bias from that, as we only have one single difference variable for each feature now.

2.1.3. Hybridizing Match & Set Features

TODO

2.1.4. Common Opponent Model

We use a method proposed by Knottenbelt [4] to construct difference variables capturing player characteristics and past performance, by using their **common opponents** so as to achieve a fair comparison. Although this technique was developed for use in a hierarchical Markov model, Spiko [8] has shown that it adapts well in the land of machine learning algorithms and modeling.

To show how we construct the features, we first look for a set of **common opponents** between the two players of a given match. Here, **common opponents** are defined as other players that the two (Player 1 and Player 2) have faced in the past on their own. Once we have all the common opponents, we find both players' performance against the common opponent (C_i), and average across all the common opponents. Lastly, we find the difference between the two players' average stats against common opponents as the selected feature. Figure 2.1 shows an example of how we could construct the **ACE** feature (difference in average ace counts) against common opponents. As shown in Figure 2.1, common opponents are labeled as C_1, C_2, \dots, C_n . For Player $i \in \{1, 2\}$, $ACE_i(C_j)$ is his/her average count of aces in all matches against common opponent C_j . We then take an average of these values to get ACE_i for Player i :

$$ACE_i = \frac{\sum_{j=0}^n ACE_i(C_j)}{n}$$

Lastly, applying the idea of difference variables, we create the feature ACE by finding the difference in performance on this feature between Player 1 and Player 2:

$$ACE = ACE_1 - ACE_2$$

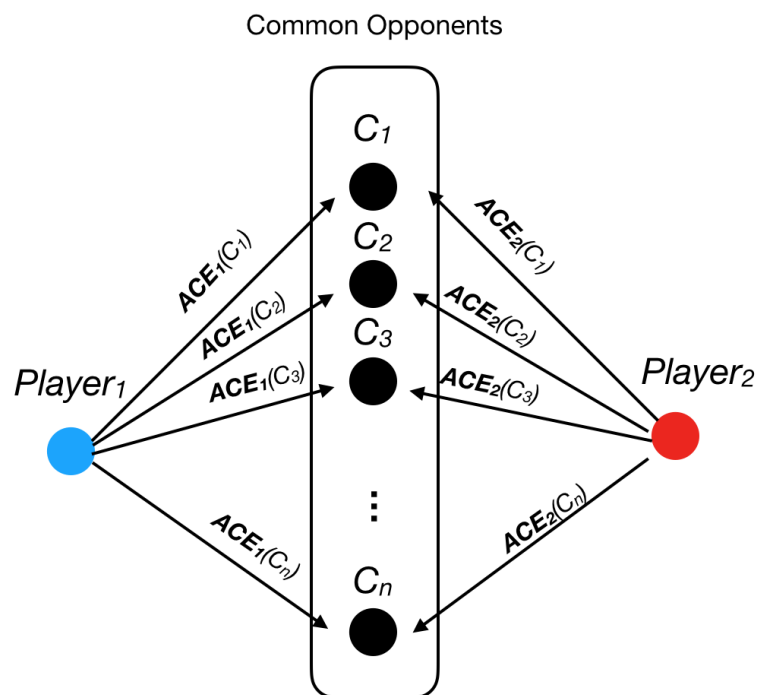
. We can apply the same technique on other features too. The limitation of this method is apparent: it only works well when players actually have common opponents have played games with them before the queried match date. Hence, we need to think about how to handle edge cases.

2.1.5. Edge Cases

We have run into 2 edge cases: (1) when at least one of the players is new; (2) when two players do not have any common opponents. When (1) happens, it is necessarily the case that it suffers from (2) no common opponents as well; but the reverse is not true.

TODO: add in what we decided to do

Figure 2.1: Constructing the ACE feature using common opponents



2.2. Data Preparation

2.2.1. Data Source

We obtain match-level data from Jeff Sackmann’s Github repositories [7] [3], with the Match Charting Project (MCP) providing set-level stats and results as well [3]. MCP is essentially a crowdsourced shot-by-shot professional tennis dataset, and it enables us to have a real-time lens into tennis matches.

2.2.2. Data Cleaning & Transformation

TODO: Things to talk about:

- Validations run on the datasets
- How to detect and remove duplicates
- Decision on dropping Davis Cup data (because team tourney might have different strategies, hence introducing biases)
- How to merge two datasets together
- (TODO for future work) Standardize features?

2.3. Feature List

We present the list of extracted features in Table 2.1. Note that all these features are difference variables unless otherwise noted.

3. Preliminary Experiments

As the team is still working on connecting the historical, match-level dataset with the real-time, set-level data sets, we ran some preliminary models fitting the two datasets separately instead. We use scikit-learn’s implementation of the various models [6].

3.1. Historical, Match-level Dataset: Logistic Regression Model

Applying a train-test split of 2:1, we trained a logistic regression model on $2/3$ of all the ATP matches from 2012. Following the labelling from Section 2.1.1, we set $label = 1$ for all the matches with Player 1 being the Winner (notation-only), and then duplicated the entire dataset but swapped Player 1 and Player 2 and set $label = 0$, as well as inverted all difference features. **TODO:** might need more explanation?

We achieved an impressive accuracy rate of 87.8% just using the historical, match-level dataset with common opponent model and difference variables.

Feature	Description
DURATION	Match duration in minutes (<i>not difference</i>)
SAME_HANDEDNESS	Tracking, on average, the frequency of facing opponents who have the same dominant hand
RANK	ATP rank
RANK_PTS	ATP rank points
HEIGHT	Player height
FSP	First serve success percentage
ACE	Average number of aces per game
DF	Average number of double faults per game
BP_FACED	Average number of break points faced per game
BP_SAVED	Average number of break points saved per game
BPP	Break point saving percentage
SVGM	Average number of serve games
W1SP	Winning on first serve percentage
W2SP	Winning on second serve percentage
SVPT	Average number of serve points per game

Table 2.1: Summary of extracted features

3.1.1. Error Analysis

A preliminary error analysis was done to explore the relative importance of all the extracted features, and the results are shown in Table 3.1. Interestingly, we note that **RANK_PTS** accounts for almost 20% of the performance.

3.2. Historical, Match-level Dataset: Naive Bayes Classifier

We also ran the data through a Gaussian-based Naive Bayes classifier, which showed an impressive 86.9% accuracy rate (misclassifying 718 points out of 5496). As expected, it was also extremely fast. The impressive performance indicates that the extracted features might indeed be conditionally independent of each other, which is something we never thought of before, and look to explore in the next steps. On a first glance, it seems to suggest that given Player i won the match, his/her performance on aces and rank point difference or double faults are indeed almost independent. We were surprised by the results, and more sanity check is needed before we put ourselves behind this.

3.3. Real-time, Set-level Dataset: Logistic Regression Model

TODO

Order	Accuracy Rate	Feature	Performance Difference
0	0.878170	SVPT	
1	0.878170	SAME_HANDEDNESS	0.00%
2	0.666483	RANK_PTS	-21.17%
3	0.657111	RANK	-0.94%
4	0.660419	HEIGHT	0.33%
5	0.659868	FSP	-0.06%
6	0.665932	DURATION	0.61%
7	0.661521	DF	-0.44%
8	0.664278	BPP	0.28%
9	0.664829	BP_SAVED	0.06%
10	0.665380	BP_FACED	0.06%
11	0.648842	ACE	-1.65%
12	0.661521	SVGM	1.27%
13	0.659868	W2SP	-0.17%

Table 3.1: Summary of extracted feature performance in the preliminary logistic regression model fitting on historical match-level data

3.4. Real-time, Set-level Dataset: SVM

TODO

4. Next Steps

The biggest next step is to connect the historical, match-level dataset with the real-time, set-level data set in order to conduct further experiments. We also plan to incorporate more sophisticated ideas on feature scaling in the context of tennis, such as **time discounting**, **surface weighting**, and the effect of fatigue and player injury. Moreover, we wanted to explore more interesting features such as **head-to-head results**, as well as to conduct **Principal Component Analysis** and **Feature Selection** to understand the inter-correlation between features, and which features are more predictive than others. Lastly, we plan to benchmark various machine learning models' performances, and to fine-tune the best them to achieve optimal performance.

5. Conclusion

The use of common opponent model and difference variables already offer high predictive power in simple machine learning models such as logistic regression and Naive Bayes, and it's exciting to see with fine-tuning and more exploration and connecting datasets what we can achieve. It's also helpful to learn that ATP rank point difference between two players alone is already predictive of their likely match outcome.

A. Contributions

A.1. Yi Zhong

A.2. Yubo Tian

A.3. Yang "Eddie" Chen

- Researched and reviewed relevant literature for applicable ideas, such as Clarke and Dyte [1] and O'Malley [2] for difference variable, and Knottenbelt [4] for the common opponent model
- Implemented the Common Opponent Model to look up players stats from common opponents
- Generated ATP Match results for labels
- Typesetted (L^AT_EX) the project milestone report

B. Appendix: Project Roadmap

- 11/13 - 11/19: Data cleaning
- 11/20 - 11/26: Preliminary experiments, building out more sophisticated ideas into models such as surface weighting, time decaying, fatigue, and injury
- 11/27 - 12/03: Continue fine-tuning models
- 12/04 - 12/10: Poster making
- 12/11 - 12/15: Poster presentation; final project write-up

References

- [1] Stephen R. Clarke and David Dyte. Using official ratings to simulate major tennis tournaments. *International Transactions in Operational Research*, 7(6):585 – 594, 2010.
- [2] O'Malley A. James. Probability formulas and statistical analysis in tennis. *Journal of Quantitative Analysis in Sports*, 4(2):1–23, April 2008.
- [3] et al. Jeff Sackmann. The tennis abstract match charting project, 2017.
- [4] William J. Knottenbelt, Demetris Spanias, and Agnieszka M. Madurska. A common-opponent stochastic model for predicting the outcome of professional tennis matches. *Computers and Mathematics with Applications*, 64(12):3820 – 3827, 2012. Theory and Practice of Stochastic Modeling.

- [5] Agnieszka M. Madurska. A set-by-set analysis method for predicting the outcome of professional singles tennis matches. MEng computing- final year project, Imperial College London, amm208@doc.ic.ac.uk, June 2012.
- [6] F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830, 2011.
- [7] Jeff Sackmann. Atp tennis rankings, results, and stats, 2017.
- [8] Michal Sipko. Machine learning for the prediction of professional tennis matches. MEng computing - final year project, Imperial College London, June 2015.